

Ancestor Sampling for Particle Gibbs

Fredrik Lindsten
Div. of Automatic Control
Linköping University
lindsten@isy.liu.se

Michael I. Jordan
Dept. of EECS and Statistics
University of California, Berkeley
jordan@cs.berkeley.edu

Thomas B. Schön
Div. of Automatic Control
Linköping University
schon@isy.liu.se

October 26, 2012

Abstract

We present a novel method in the family of particle MCMC methods that we refer to as *particle Gibbs with ancestor sampling* (PG-AS). Similarly to the existing *PG with backward simulation* (PG-BS) procedure, we use backward sampling to (considerably) improve the mixing of the PG kernel. Instead of using separate forward and backward sweeps as in PG-BS, however, we achieve the same effect in a single forward sweep. We apply the PG-AS framework to the challenging class of non-Markovian state-space models. We develop a truncation strategy of these models that is applicable in principle to any backward-simulation-based method, but which is particularly well suited to the PG-AS framework. In particular, as we show in a simulation study, PG-AS can yield an order-of-magnitude improved accuracy relative to PG-BS due to its robustness to the truncation error. Several application examples are discussed, including Rao-Blackwellized particle smoothing and inference in degenerate state-space models. This report is a slightly extended version of the paper [1].

1 Introduction

State-space models (SSMs) are widely used to model time series and dynamical systems. The strong assumptions of linearity and Gaussianity that were originally invoked in state-space inference have been weakened by two decades of research on sequential Monte Carlo (SMC) and Markov chain Monte Carlo (MCMC). These Monte Carlo methods have not, however, led to substantial weakening of a further strong assumption, that of Markovianity. It remains a

major challenge to develop inference algorithms for non-Markovian SSMs:

$$x_{t+1} \sim f(x_{t+1} \mid \theta, x_{1:t}), \quad y_t \sim g(y_t \mid \theta, x_{1:t}), \quad (1)$$

where $\theta \in \Theta$ is a static parameter with prior density $p(\theta)$, x_t is the latent state and y_t is the observation at time t , respectively. Models of this form arise in many different application scenarios, either from direct modeling or via a transformation or marginalization of a larger model. We provide several examples in Section 5.

To tackle the challenging problem of inference for non-Markovian SSMs, we work within the framework of particle MCMC (PMCMC), a family of inferential methods introduced in [2]. The basic idea in PMCMC is to use SMC to construct a proposal kernel for an MCMC sampler. Assume that we observe a sequence of measurements $y_{1:T}$. We are interested in finding the density $p(x_{1:T}, \theta \mid y_{1:T})$, i.e., the joint posterior density of the state sequence and the parameter. In an idealized Gibbs sampler we would target this density by sampling as follows: (i) Draw $\theta^* \mid x_{1:T} \sim p(\theta \mid x_{1:T}, y_{1:T})$; (ii) Draw $x_{1:T}^* \mid \theta^* \sim p(x_{1:T} \mid \theta^*, y_{1:T})$. The first step of this procedure can be carried out exactly if conjugate priors are used. For non-conjugate models, one option is to replace Step (i) with a Metropolis-Hastings step. However, Step (ii)—sampling from the joint smoothing density $p(x_{1:T} \mid \theta, y_{1:T})$ —is in most cases very difficult. In PMCMC, this is addressed by instead sampling a particle trajectory $x_{1:T}^*$ based on an SMC approximation of the joint smoothing density. More precisely, we run an SMC sampler targeting $p(x_{1:T} \mid \theta^*, y_{1:T})$. We then sample one of the particles at the final time T , according to their importance weights, and trace the ancestral lineage of this particle to obtain the trajectory $x_{1:T}^*$. This overall procedure is referred to as *particle Gibbs* (PG).

The flexibility provided by the use of SMC as a proposal mechanism for MCMC seems promising for tackling inference in non-Markovian models. To exploit this flexibility we must address a drawback of PG in the high-dimensional setting, which is that the mixing of the PG kernel can be very poor when there is path degeneracy in the SMC sampler [3, 4]. This problem has been addressed in the generic setting of SSMs by adding a backward simulation step to the PG sampler, yielding a method denoted *PG with backward simulation* (PG-BS). It has been found that this considerably improves mixing, making the method much more robust to a small number of particles as well as larger data records [3, 4].

Unfortunately, however, the application of backward simulation is problematic for non-Markovian models. The reason is that we need to consider full state trajectories during the backward simulation pass, leading to $O(T^2)$ computational complexity (see Section 4 for details). To address this issue, we develop a novel PMCMC method which we refer to as *particle Gibbs with ancestor sampling* (PG-AS) that achieves the effect of backward sampling without an explicit backward pass. As part of our development, we also develop a truncation method geared to non-Markovian models. This method is a generic method that is also applicable to PG-BS, but, as we show in a simulation study in Section 6, the effect of the truncation error is much less severe for PG-AS than for

PG-BS. Indeed, we obtain up to an order of magnitude increase in accuracy in using PG-AS when compared to PG-BS in this study.

Since we assume that it is straightforward to sample the parameter θ of the idealized Gibbs sampler, we will not explicitly include sampling of θ in the subsequent sections to simplify our presentation.

This report is a slightly extended version of the paper [1].

2 Sequential Monte Carlo

We first review the standard auxiliary SMC sampler, see e.g. [5, 6]. Let $\gamma_t(x_{1:t})$ for $t = 1, \dots, T$ be a sequence of unnormalized densities on \mathbf{X}^t , which we assume can be evaluated pointwise in linear time. Let $\bar{\gamma}_t(x_{1:t})$ be the corresponding normalized probability densities. For an SSM we would typically have $\bar{\gamma}_t(x_{1:t}) = p(x_{1:t} \mid y_{1:t})$ and $\gamma_t(x_{1:t}) = p(x_{1:t}, y_{1:t})$. Assume that $\{x_{1:t-1}^m, w_{t-1}^m\}_{m=1}^N$ is a weighted particle system targeting $\bar{\gamma}_{t-1}(x_{1:t-1})$. This particle system is propagated to time t by sampling independently from a proposal kernel,

$$M_t(a_t, x_t) = \frac{w_{t-1}^{a_t} \nu_{t-1}^{a_t}}{\sum_l w_{t-1}^l \nu_{t-1}^l} R_t(x_t \mid x_{1:t-1}^{a_t}). \quad (2)$$

In this formulation, the resampling step is implicit and corresponds to sampling the ancestor indices a_t . Note that a_t^m is the index of the ancestor particle of x_t^m . When we write $x_{1:t}^m$ we refer to the ancestral path of x_t^m . The factors $\nu_t^m = \nu_t(x_{1:t}^m)$, known as adjustment multiplier weights, are used in the auxiliary SMC sampler to increase the probability of sampling ancestors that better can describe the current observation [6]. The particles are then weighted according to $w_t^m = W_t(x_{1:t}^m)$, where the weight function is given by

$$W_t(x_{1:t}) = \frac{\gamma_t(x_{1:t})}{\gamma_{t-1}(x_{1:t-1}) \nu_{t-1}(x_{1:t-1}) R_t(x_t \mid x_{1:t-1})}, \quad (3)$$

for $t \geq 2$. The procedure is initiated by sampling from a proposal density $x_1^m \sim R_1(x_1)$ and assigning importance weights $w_1^m = W_1(x_1^m)$ with $W_1(x_1) = \gamma_1(x_1)/R_1(x_1)$. In PMCMC it is instructive to view this sampling procedure as a way of generating a single sample from the density

$$\psi(\mathbf{x}_{1:T}, \mathbf{a}_{2:T}) \triangleq \prod_{m=1}^N R_1(x_1^m) \prod_{t=2}^T \prod_{m=1}^N M_t(a_t^m, x_t^m) \quad (4)$$

on the space $\mathbf{X}^{NT} \times \{1, \dots, N\}^{N(T-1)}$. Here we have introduced the boldface notation $\mathbf{x}_t = \{x_t^1, \dots, x_t^N\}$ and similarly for the ancestor indices.

3 Particle Gibbs with ancestor sampling

PMCMC methods is a class of MCMC samplers in which SMC is used to construct proposal kernels [2]. The validity of these methods can be assessed by

viewing them as MCMC samplers on an extended state space in which all the random variables generated by the SMC sampler are seen as auxiliary variables. The target density on this extended space is given by

$$\phi(\mathbf{x}_{1:T}, \mathbf{a}_{2:T}, k) \triangleq \frac{\bar{\gamma}_T(x_{1:T}^k)}{N^T} \frac{\psi(\mathbf{x}_{1:T}, \mathbf{a}_{2:T})}{R_1(x_1^{b_1}) \prod_{t=2}^T M_t(a_t^{b_t}, x_t^{b_t})}. \quad (5)$$

By construction, this density admits $\bar{\gamma}_T(x_{1:T}^k)$ as a marginal, and can thus be used as a surrogate for the original target density $\bar{\gamma}_T$ [2]. Here k is a variable indexing one of the particles at the final time point and $b_{1:T}$ corresponds to the ancestral path of this particle: $x_{1:T}^k = x_{1:T}^{b_{1:T}} = \{x_1^{b_1}, \dots, x_T^{b_T}\}$. These indices are given recursively from the ancestor indices by $b_T = k$ and $b_t = a_{t+1}^{b_{t+1}}$. The PG sampler [2] is a Gibbs sampler targeting ϕ using the following sweep (note that $b_{1:T} = \{a_{2:T}^{b_{2:T}}, b_T\}$),

1. Draw $\mathbf{x}_{1:T}^{*, -b_{1:T}}, \mathbf{a}_{2:T}^{*, -b_{2:T}} \sim \phi(\mathbf{x}_{1:T}^{-b_{1:T}}, \mathbf{a}_{2:T}^{-b_{2:T}} \mid x_{1:T}^{b_{1:T}}, b_{1:T})$.
2. Draw $k^* \sim \phi(k \mid \mathbf{x}_{1:T}^{*, -b_{1:T}}, \mathbf{a}_{2:T}^{*, -b_{2:T}}, x_{1:T}^{b_{1:T}}, a_{2:T}^{b_{2:T}})$.

Here we have introduced the notation $\mathbf{x}_t^{-m} = \{x_t^1, \dots, x_t^{m-1}, x_t^{m+1}, \dots, x_t^N\}$, $\mathbf{x}_{1:T}^{-b_{1:T}} = \{\mathbf{x}_1^{-b_1}, \dots, \mathbf{x}_T^{-b_T}\}$ and similarly for the ancestor indices. In [2], a sequential procedure for sampling from the conditional density appearing in Step 1 is given. This method is known as *conditional SMC* (CSMC). It takes the form of an SMC sampler in which we condition on the event that a prespecified path $x_{1:T}^{b_{1:T}} = x'_{1:T}$, with indices $b_{1:T}$, is maintained throughout the sampler (see Algorithm 1 for a related procedure). Furthermore, the conditional distribution appearing in Step 2 of the PG sampler is shown to be proportional to w_T^k , and it can thus straightforwardly be sampled from.

Note that we never sample new values for the variables $\{x_{1:T}^{b_{1:T}}, b_{1:T-1}\}$ in this sweep. Hence, the PG sampler is an “incomplete” Gibbs sampler, since it does not loop over all the variables of the model. It still holds that the PG sampler is ergodic, which intuitively can be explained by the fact that the collection of variables that is left out is chosen randomly at each iteration. However, it has been observed that the PG sampler can have very poor mixing, especially when N is small and/or T is large [3, 4]. The reason for this poor mixing is that the SMC path degeneracy causes the collections of variables that are left out at any two consecutive iterations to be strongly dependent.

We now turn to our new procedure, PG-AS, which aims to address this fundamental issue. Our idea is to sample new values for the *ancestor indices* $b_{1:T-1}$ as part of the CSMC procedure¹. By adding these variables to the Gibbs sweep, we can considerably improve the mixing of the PG kernel. The CSMC method is a sequential procedure to sample from $\phi(\mathbf{x}_{1:T}^{-b_{1:T}}, \mathbf{a}_{2:T}^{-b_{2:T}} \mid x_{1:T}^{b_{1:T}}, b_{1:T})$ by sampling according to $\{\mathbf{x}_t^{*, -b_t}, \mathbf{a}_t^{*, -b_t}\} \sim \phi(\mathbf{x}_t^{-b_t}, \mathbf{a}_t^{-b_t} \mid \mathbf{x}_{1:t-1}^{*, -b_{1:t-1}}, \mathbf{a}_{2:t-1}^{*, -b_{2:t-1}}, x_{1:T}^{b_{1:T}}, b_{1:T})$,

¹Ideally, we would like to include the variables $x_{1:T}^{b_{1:T}}$ as well, but this is in general not possible since it would be similar to sampling from the original target density (which we assume is infeasible).

for $t = 1, \dots, T$. After having sampled these variables at time t , we add a step in which we generate a new value for $b_{t-1}(=a_t^{b_t})$, resulting in the following sweep:

1'. (CSMC with ancestor sampling) For $t = 1, \dots, T$, draw

$$\begin{aligned} \mathbf{x}_t^{*, -b_t}, \mathbf{a}_t^{*, -b_t} &\sim \phi(\mathbf{x}_t^{-b_t}, \mathbf{a}_t^{-b_t} \mid \mathbf{x}_{1:t-1}^{*, -b_{1:t-1}}, \mathbf{a}_{2:t-1}^*, x_{1:T}^{b_{1:T}}, b_{t-1:T}), \\ (a_t^{*, b_t} =) b_{t-1}^* &\sim \phi(b_{t-1} \mid \mathbf{x}_{1:t-1}^{*, -b_{1:t-1}}, \mathbf{a}_{2:t-1}^*, x_{1:T}^{b_{1:T}}, b_{t:T}). \end{aligned}$$

2'. Draw $(k^* =) b_T^* \sim \phi(b_T \mid \mathbf{x}_{1:T}^{*, -b_{1:T}}, \mathbf{a}_{2:T}^*, x_{1:T}^{b_{1:T}})$.

It can be verified that this corresponds to a partially collapsed Gibbs sampler [7] and will thus leave ϕ invariant. To determine the conditional densities from which the ancestor indices are drawn, consider the following factorization, following directly from (3),

$$\begin{aligned} \gamma_t(x_{1:t}) &= W_t(x_{1:t}) \nu_{t-1}(x_{1:t-1}) R_t(x_t \mid x_{1:t-1}) \gamma_{t-1}(x_{1:t-1}) \\ \Rightarrow \gamma_t(x_{1:t}^{b_t}) &= w_t^{b_t} \frac{\sum_l w_{t-1}^l \nu_{t-1}^l}{w_{t-1}^{b_{t-1}}} \frac{w_{t-1}^{b_{t-1}} \nu_{t-1}^{b_{t-1}}}{\sum_l w_{t-1}^l \nu_{t-1}^l} R_t(x_t^{b_t} \mid x_{1:t-1}^{b_{t-1}}) \gamma_{t-1}(x_{1:t-1}^{b_{t-1}}) \\ &= \dots = w_t^{b_t} \left(\prod_{s=1}^{t-1} \sum_l w_s^l \nu_s^l \right) R_1(x_1^{b_1}) \prod_{s=2}^t M_t(a_s^{b_s}, x_s^{b_s}). \end{aligned} \quad (6)$$

Furthermore, we have

$$\begin{aligned} \phi(b_t \mid \mathbf{x}_{1:t}, \mathbf{a}_{2:t}, x_{t+1:T}^{b_{t+1:T}}, b_{t+1:T}) &\propto \phi(\mathbf{x}_{1:t}, \mathbf{a}_{2:t}, x_{t+1:T}^{b_{t+1:T}}, b_{t:T}) \\ &\propto \frac{\gamma_T(x_{1:T}^k) \psi(\mathbf{x}_{1:t}, \mathbf{a}_{2:t})}{R_1(x_1^{b_1}) \prod_{s=2}^t M_s(a_s^{b_s}, x_s^{b_s})} \propto \frac{\gamma_t(x_{1:t}^{b_t})}{\gamma_t(x_{1:t}^{b_t})} \frac{\gamma_T(x_{1:T}^k)}{R_1(x_1^{b_1}) \prod_{s=2}^t M_s(a_s^{b_s}, x_s^{b_s})}. \end{aligned} \quad (7)$$

By plugging (6) into the numerator we get,

$$\phi(b_t \mid \mathbf{x}_{1:t}, \mathbf{a}_{2:t}, x_{t+1:T}^{b_{t+1:T}}, b_{t+1:T}) \propto w_t^{b_t} \frac{\gamma_T(x_{1:T}^k)}{\gamma_t(x_{1:t}^{b_t})}. \quad (8)$$

Hence, to sample a new ancestor index for the conditioned path at time $t+1$, we proceed as follows. Given $x'_{t+1:T} (= x_{t+1:T}^{b_{t+1:T}})$ we compute the backward sampling weights,

$$w_{t|T}^m = w_t^m \frac{\gamma_T(\{x_{1:t}^m, x'_{t+1:T}\})}{\gamma_t(x_{1:t}^m)}, \quad (9)$$

for $m = 1, \dots, N$. We then set $b_t = m$ with probability proportional to $w_{t|T}^m$.

It follows that the proposed CSMC with ancestor sampling (Step 1'), conditioned on $\{x'_{1:T}, b_{1:T}\}$, can be realized as in Algorithm 1. The difference between this algorithm and the CSMC sampler derived in [2] lies in the ancestor sampling step 2(b) (where instead, they set $a_t^{b_t} = b_{t+1}$). By introducing the ancestor

Algorithm 1 CSMC with ancestor sampling, conditioned on $\{x'_{1:T}, b_{1:T}\}$

1. **Initialize** ($t = 1$):
 - (a) Draw $x_1^m \sim R_1(x_1)$ for $m \neq b_1$ and set $x_1^{b_1} = x'_1$.
 - (b) Set $w_1^m = W_1(x_1^m)$ for $m = 1, \dots, N$.
 2. **for** $t = 2, \dots, T$:
 - (a) Draw $\{a_t^m, x_t^m\} \sim M_t(a_t, x_t)$ for $m \neq b_t$ and set $x_t^{b_t} = x'_t$.
 - (b) Draw $a_t^{b_t}$ with $P(a_t^{b_t} = m) \propto w_{t-1|T}^m$.
 - (c) Set $x_{1:t}^m = \{x_{1:t-1}^m, x_t^m\}$ and $w_t^m = W_t(x_{1:t}^m)$ for $m = 1, \dots, N$.
-

sampling, we break the strong dependence between the generated particle trajectories and the path on which we condition. We call the resulting method, defined by Steps 1' and 2' above, *PG with ancestor sampling* (PG-AS).

The idea of including the variables $b_{1:T-1}$ in the PG sampler has previously been suggested by Whiteley [8] and further explored in [3,4]. This previous work, however, accomplishes this with an explicit backward simulation pass, which, as we discuss in the following section, is problematic for our applications to non-Markovian SSMs. In the PG-AS sampler, instead of requiring distinct forward and backward sequences of Gibbs steps as in PG with backward simulation (PG-BS), we obtain a similar effect via a single forward sweep.

4 Truncation for non-Markovian models

We return to the problem of inference in non-Markovian SSMs of the form shown in (1). To employ backward sampling, we need to evaluate the ratio

$$\frac{\gamma_T(x_{1:T})}{\gamma_t(x_{1:t})} = \frac{p(x_{1:T}, y_{1:T})}{p(x_{1:t}, y_{1:t})} = \prod_{s=t+1}^T g(y_s | x_{1:s}) f(x_s | x_{1:s-1}). \quad (10)$$

In general, the computational cost of computing the backward sampling weights will thus be $O(T)$. This implies that the cost of generating a full backward trajectory is $O(T^2)$. It is therefore computationally prohibitive to employ backward simulation type of particle smoothers, as well as the PG samplers discussed above, for general non-Markovian models.

To make progress, we consider non-Markovian models in which there is a decay in the influence of the past on the present, akin to that in Markovian models but without the strong Markovian assumption. Hence, it is possible to obtain a useful approximation when the product in (10) is truncated to a

smaller number of factors, say p . We then replace (9) with the approximation,

$$\tilde{w}_{t|T}^{p,m} = w_t^m \frac{\gamma_{t+p}(\{x_{1:t}^m, x'_{t+1:t+p}\})}{\gamma_t(x_{1:t}^m)}. \quad (11)$$

The following proposition formalizes our assumption.

Proposition 1. *Let P and \tilde{P}_p be the probability distributions on $\{1, \dots, N\}$, defined by the backward sampling weight (9) and the truncated backward sampling weights (11), respectively. Let $h_s(k) = g(y_{t+s} \mid x_{1:t}^k, x'_{t+1:t+s})f(x'_{t+s} \mid x_{1:t}^k, x'_{t+1:t-s})$ and assume that $\max_{k,l} (h_s(k)/h_s(l) - 1) \leq A \exp(-cs)$, for some constants A and $c > 0$. Then, $D_{\text{KLD}}(P \parallel \tilde{P}_p) \leq C \exp(-cp)$ for some constant C , where D_{KLD} is the Kullback-Leibler divergence (KLD).*

Proof. See Appendix A. □

From (11), we see that we can compute the backward weights in constant time under the truncation within the PG-AS framework. The resulting approximation can be quite useful; indeed, in our experiments we have seen that even $p = 1$ can lead to very accurate inferential results. In general, however, it will not be known a priori how to set the truncation level p for any given problem. To address this problem, we propose to use an adaption of the truncation level. Since the approximative weights (11) can be evaluated sequentially, the idea is to start with $p = 1$ and then increase p until the weights have, in some sense, converged. In particular, in our experimental work, we have used the following simple approach.

Let P_p be the discrete probability measure defined by (11). Let $\varepsilon_p = D_{\text{TV}}(\tilde{P}_p, \tilde{P}_{p-1})$ be the total variation (TV) distance between the distributions for two consecutive truncation levels. We then compute the exponentially decaying moving average of the sequence ε_p , with forgetting factor $\gamma \in [0, 1]$, and stop when this falls below some threshold $\tau \in [0, 1]$. This adaption scheme removes the requirement to specify p directly, but instead introduces the design parameters γ and τ . However, these parameters are much easier to reason about – a small value for γ gives a rapid response to changes in ε_p whereas a large value gives a more conservative stopping rule, improving the accuracy of the approximation at the cost of higher computational complexity. A similar trade off holds for the threshold τ as well. Most importantly, we have found that the same values for γ and τ can be used for a wide range of models, with very different mixing properties.

To illustrate the effect of the adaption rule, and how the distribution \tilde{P}_p typically evolves as we increase p , we provide two examples in Figure 1. These examples are taken from the simulation study provided in Section 6.2. Note that the untruncated distribution P is given for the maximal value of p , i.e., furthest to the right in the figures. By using the adaptive truncation, we can stop the evaluation of the weights at a much earlier stage, and still obtain an accurate approximation of P .

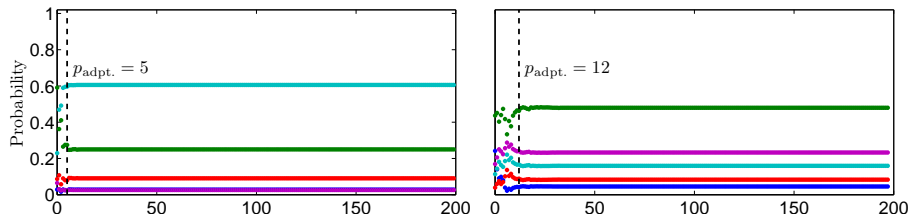


Figure 1: Probability under \tilde{P}_p as a function of the truncation level p for two different systems; one 5 dimensional (left) and one 20 dimensional (right). The $N = 5$ dotted lines correspond to $\tilde{P}_p(m)$ for $m \in \{1, \dots, N\}$, respectively (N.B. two of the lines overlap in the left figure). The dashed vertical lines show the value of the truncation level $p_{\text{adpt.}}$, resulting from the adaption scheme with $\gamma = 0.1$ and $\tau = 10^{-2}$. See Section 6.2 for details on the experiments.

5 Application areas

In this section we present examples of problem classes involving non-Markovian SSMs for which the proposed PG-AS sampler can be applied. Numerical illustrations are provided in Section 6.

5.1 Rao-Blackwellized particle smoothing

One popular approach to increase the efficiency of SMC samplers for SSMs is to marginalize over one component of the state, and apply an SMC sampler in the lower-dimensional marginal space. This leads to what is known as the Rao-Blackwellized particle filter (RBPF) [9–11]. The same approach has also been applied to state smoothing [12, 13], but it turns out that Rao-Blackwellization is less straightforward in this case, since the marginal state-process will be non-Markovian. As an example, a mixed linear/nonlinear Gaussian SSM (see, e.g., [11]) with “nonlinear state” x_t and “conditionally linear state” z_t , can be reduced to

$$x_t \sim p(x_t \mid x_{1:t-1}, y_{1:t-1}), \quad y_t \sim p(y_t \mid x_{1:t}, y_{1:t-1}). \quad (12)$$

These conditional densities are Gaussian and can be evaluated for any fixed marginal state trajectory $x_{1:t-1}$ by running a conditional Kalman filter to marginalize the z_t -process.

In order to apply a backward-simulation-based method (e.g., a particle smoother) for this model, we need to evaluate the backward sampling weights (9). In a straightforward implementation, we thus need to run N Kalman filters for $T - t$ time steps, for each $t = 1, \dots, T - 1$. The computational complexity of this calculation can be reduced by employing the truncation proposed in Section 4².

²For the specific problem of Rao-Blackwellized smoothing in conditionally Gaussian models, a backward simulator which can be implemented in $O(T)$ computational complexity has recently been proposed in [12]. This is based on the idea of propagating information backward in time as the backward samples are generated.

5.2 Particle smoothing for degenerate state-space models

Many dynamical systems are most naturally modelled as degenerate in the sense that the transition kernel of the state process does not admit any dominating measure. For instance, consider a nonlinear system with additive noise of the form,

$$\xi_t = f(\xi_{t-1}) + G\omega_{t-1}, \quad y_t = g(\xi_t) + e_t, \quad (13)$$

where G is a tall matrix, and consequently $\text{rank}(G) < \dim(\xi_t)$. That is, the process noise covariance matrix is singular. SMC samplers can straightforwardly be applied to this type of models, but it is more problematic to address the smoothing problem using particle methods. The reason is that the backward kernel also will be degenerate and it cannot be approximated in a natural way by the forward filter particles, as is normally done in backward-simulation-based particle smoothers.

A possible remedy for this issue is to recast the degenerate SSM as a non-Markovian model in a lower-dimensional space. Let $G = U [\Sigma \ 0]^\top V^\top$ with unitary U and V be a singular value decomposition of G and let,

$$\begin{bmatrix} x_t \\ z_t \end{bmatrix} \triangleq U^\top \xi_t = U^\top f(UU^\top \xi_{t-1}) + \begin{bmatrix} \Sigma V^\top \omega_{t-1} \\ 0 \end{bmatrix}. \quad (14)$$

For simplicity we assume that z_1 is known. If this is not the case, it can be included in the system state or seen as a static parameter of the model. Hence, the sequence $z_{1:t}$ is $\sigma(x_{1:t-1})$ -measurable and we can write $z_t = z_t(x_{1:t-1})$. With $v_t \triangleq \Sigma V^\top \omega_t$ and by appropriate definitions of the functions f_x and h , the model (13) can thus be rewritten as,

$$x_t = f_x(x_{1:t-1}) + v_{t-1}, \quad y_t = h(x_{1:t}) + e_t, \quad (15)$$

which is a non-degenerate, non-Markovian SSM. By exploiting the truncation proposed in Section 4 we can thus apply PG-AS to do inference in this model. In fact, this is nothing but another application of Rao-Blackwellization as discussed in Section 5.1, where the z_t -state is conditionally deterministic and thus trivially marginalizable.

5.3 Additional problem classes

There are many more problem classes in which non-Markovian models arise and in which backward-simulation-based methods can be of interest. For instance, the Dirichlet process mixture model (DPMM, see, e.g., [14]) is a popular nonparametric Bayesian model for mixtures with an unknown number of components. Using a Polya urn representation, the mixture labels are given by a non-Markovian stochastic process, and the DPMM can thus be seen as a non-Markovian SSM. SMC has previously been used for inference in DPMMs [15,16]. An interesting venue for future work is to use the PG-AS sampler for these models. A second example in Bayesian nonparametrics is Gaussian

process (GP) regression and classification (see, e.g., [17]). The sample path of the GP can be seen as the state-process in a non-Markovian SSM. We can thus employ PMCMC, and in particular PG-AS, to address these inference problems.

An application in genetics, for which SMC has been successfully applied, is reconstruction of phylogenetic trees [18]. A phylogenetic tree is a binary tree with observation at the leaf nodes. SMC is used to construct the tree in a bottom up fashion. A similar approach has also been used for Bayesian agglomerative clustering, in which SMC is used to construct a binary clustering tree based on Kingman’s coalescent [19]. The generative models for the trees used in [18, 19] are in fact Markovian, but the observations give rise to a conditional dependence which destroys the Markov property. To employ backward simulation to these models, we are thus faced with problems of a similar nature as those discussed in Section 4.

6 Numerical evaluation

This section contains a numerical evaluation of the proposed method. First, we consider linear Gaussian systems, which is instructive since the exact smoothing density then is available, e.g., by running a modified Bryson-Frazier (MBF) smoother [20]. Second, we apply the proposed method for joint state and parameter inference in a target tracking scenario.

6.1 RBPS: Linear Gaussian state-space model

As a first example, we consider Rao-Blackwellized particle smoothing (RBPS) in a single-output 4th-order linear Gaussian SSM. The system has poles in -0.65 , -0.12 and $0.22 \pm 0.10i$ and is excited by white Gaussian noise with variance $0.1I_4$. The scalar output y_t is observed in white Gaussian noise with variance 0.1 . We generate $T = 100$ samples from the system and run PG-AS and PG-BS, marginalizing three out of the four states using an RBPF, i.e., $\dim(x_t) = 1$. Both methods are run for $R = 10000$ iterations using $N = 5$ particles. The truncation level is set to $p = 1$, leading to a coarse approximation. The total computational complexity for each sampler is $O(RNTp)$. We discard the first 1000 iterations and then compute running means of the state trajectory $x_{1:T}$. From these, we then compute the running root mean squared errors (RMSEs) ϵ_r *relative to the true posterior means* (computed with an MBF smoother). Hence, if no approximation would have been made, we would expect $\epsilon_r \rightarrow 0$, so any static error can be seen as the effect of the truncation. The results for five independent runs from both PG samplers are shown in Figure 2. First, we note that both methods give accurate results. Still, the error for PG-AS is close to an order of magnitude less than for PG-BS. Furthermore, it appears as if the error for PG-AS would decrease further, given more iterations, suggesting that the bias caused by the truncation is dominated by the Monte Carlo variance, even after $R = 10000$ iterations.

For further comparison, we also run an untruncated forward filter/backward

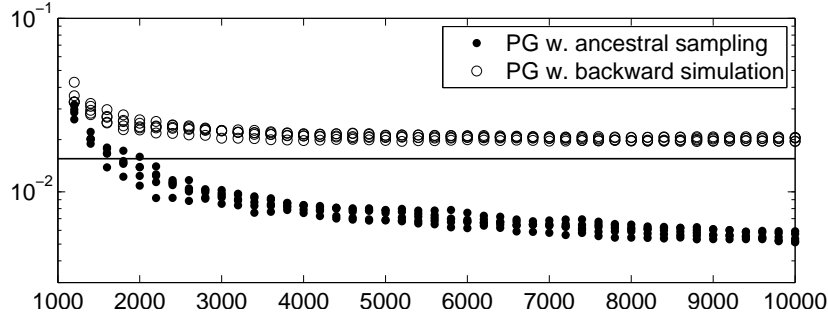


Figure 2: Rao-Blackwellized state smoothing using PG. Running RMSEs for five independent runs of PG-AS (\bullet) and PG-BS (\circ), respectively. The truncation level is set to $p = 1$. The solid line corresponds to a run of an untruncated FF-BS.

simulator (FF-BS) particle smoother [21], using $N = 5000$ forward filter particles and $M = 500$ backward trajectories (with a computational complexity of $O(NMT^2)$). The resulting RMSE value is shown as a solid line in Figure 2. These results suggest that PMCMC samplers, such as the PG-AS, indeed can be serious competitors to more “standard” particle smoothers. Even with $p = 1$, PG-AS outperforms FF-BS in terms of accuracy and, due to the fact that the ancestor sampling allows us to use as few as $N = 5$ particles at each iteration, at a lower computational cost.

6.2 Random linear Gaussian systems with rank deficient process noise covariances

To see how the PG samplers are affected by the choice of truncation level p and by the mixing properties of the system, we evaluate them on random linear Gaussian SSMs of different orders. We generate 150 random systems, using the MATLAB function `drss` from the Control Systems Toolbox, with model orders 2, 5 and 20 (50 systems for each model order). The number of outputs are taken as 1, 2 and 4 for the different model orders, respectively. The systems are then simulated for $T = 200$ time steps, driven by Gaussian process noise entering only on the first state component. Hence, the rank of the process noise covariance is 1 for all systems. The process noise and measurement noise variances are both set to 0.1.

We run the PG-AS and PG-BS samplers for 10000 iterations using $N = 5$ particles. We consider different fixed truncation levels, ($p = 1, 2$ and 3 for 2nd order systems and $p = 1, 5$ and 10 for 5th and 20th order systems), as well as an adaptive level with $\gamma = 0.1$ and $\tau = 10^{-2}$. Again, we compute running posterior means (discarding 1000 samples) and RMSE values relative the true posterior mean. Box plots are shown in Figure 3. Since the process noise only enters on one of the state components, the mixing tends to deteriorate as we increase the

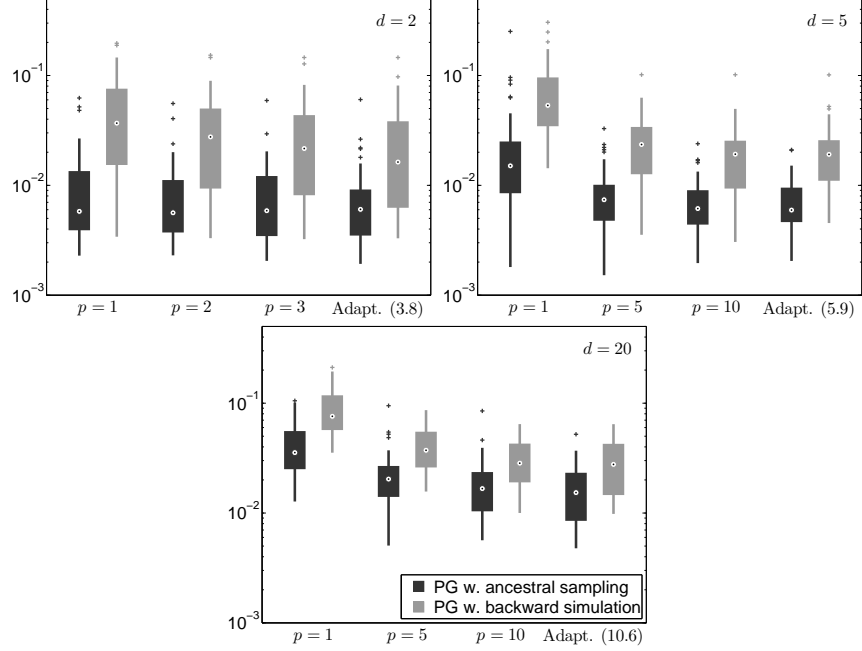


Figure 3: Box plots of the RMSE errors for PG-AS (black) and PG-BS (gray), for 150 random systems of different dimensions d (upper left, $d = 2$; upper right, $d = 5$; bottom, $d = 20$). Different values for the truncation level p are considered. The rightmost boxes correspond to an adaptive threshold and the values in parentheses are the average over all systems and MCMC iterations (the same for both methods). The dots within the boxes show the median errors.

model order. Figure 1 shows how the probability distributions on $\{1, \dots, N\}$ change as we increase the truncation level, in two representative cases for a 5th and a 20th order system, respectively. By using an adapted level, we can obtain accurate results for systems of different dimensions, without having to change any settings between the runs.

6.3 Range-bearing tracking in model with rank deficient process noise covariance

Target tracking is an area in which SMC methods have been applied with great success, see e.g. [22–24]. Tracking is most commonly seen as an online filtering problem, though in certain scenarios it might be beneficial to instead view it as a smoothing problem. For instance, if a target tracker in a surveillance system detects some abnormal behaviour, it can be interesting to apply a smoother to obtain refined estimates of the target’s position prior to the detection.

Here, we consider smoothing in a range-bearing target tracking scenario.

The system state consists of the target's position and velocity in two dimensions, $\xi_t = (p_t^x \ p_t^y \ v_t^x \ v_t^y)^\top$. We use a coordinated turn (CT) model, which is a standard model for a manoeuvring target (see e.g. [23]),

$$\begin{pmatrix} p_t^x \\ p_t^y \\ v_t^x \\ v_t^y \end{pmatrix}^\top = \underbrace{\begin{pmatrix} p_{t-1}^x + \frac{\sin(T_s \Phi_{t-1})}{\Phi_{t-1}} v_{t-1}^x - \frac{1 - \cos(T_s \Phi_{t-1})}{\Phi_{t-1}} v_{t-1}^y \\ p_{t-1}^y + \frac{1 - \cos(T_s \Phi_{t-1})}{\Phi_{t-1}} v_{t-1}^x + \frac{\sin(T_s \Phi_{t-1})}{\Phi_{t-1}} v_{t-1}^y \\ \cos(T_s \Phi_{t-1}) v_{t-1}^x - \sin(T_s \Phi_{t-1}) v_{t-1}^y \\ \sin(T_s \Phi_{t-1}) v_{t-1}^x + \cos(T_s \Phi_{t-1}) v_{t-1}^y \end{pmatrix}}_{=f_\theta(\xi_{t-1})} + \underbrace{\begin{pmatrix} \frac{T_s^2}{2} & 0 \\ 0 & \frac{T_s^2}{2} \\ T_s & 0 \\ 0 & T_s \end{pmatrix}}_{=G} \omega_{t-1}. \quad (16a)$$

The turn rate is given by

$$\Phi_t = \frac{\theta}{\sqrt{(v_t^x)^2 + (v_t^y)^2}}, \quad (16b)$$

which depends nonlinearly on the system state. The parameter θ is the manoeuvre acceleration, which we assume is fixed but unknown. This is done to illustrate the fact that PG-AS straightforwardly can be used for joint parameter and state inference, as pointed out in Section 1. The system is assumed to be affected by a random acceleration $\omega_t \sim \mathcal{N}(0, Q)$ (the process noise), here with $Q = 10I_2$. This is a common assumption for many models used in target tracking. The matrix G arises from a time discretization of a continuous time model, where $T_s = 0.1$ is the sampling time. The initial state of the system is given a Gaussian prior, $\mathcal{N}\left(\begin{pmatrix} 500 & 500 & 0 & 0 \end{pmatrix}^\top, \text{diag}\left(\begin{pmatrix} 20 & 20 & 5 & 5 \end{pmatrix}^\top\right)\right)$.

We assume that the range and bearing of the target can be observed, so that the measurements are given by,

$$y_t = \begin{pmatrix} \sqrt{(p_t^x)^2 + (p_t^y)^2} \\ \arctan(p_t^x/p_t^y) \end{pmatrix} + e_t, \quad e_t \sim \mathcal{N}\left(0, \begin{pmatrix} 50 & 0 \\ 0 & 10^{-4} \end{pmatrix}\right). \quad (17)$$

This choice of measurement noise covariance corresponds to an accurate bearing measurement, but an uninformative range measurement. Such a measurement could for instance arise in visual tracking, where the range is estimated based on the size of the target.

We initialize the system as $\xi_1 = (490 \ 490 \ 0 \ 5)^\top$ and simulate it for $T = 200$ time steps. The true target trajectory is shown in Figure 4. Note that the process noise covariance GQG^\top is singular, which implies that care needs to be taken when designing a smoothing algorithm for this model. Here, we apply a linear state transformation, as suggested in Section 5.2, to reduce the model to a lower-dimensional state-space. With $G = U[\Sigma \ 0]^\top V^\top$ we define $[x_t^\top \ z_t^\top]^\top = U^\top \xi_t$. We then employ the PG-AS sampler for joint parameter and state inference, by targeting the density $p(\theta, z_1, x_{1:T} \mid y_{1:T})$. We apply a Metropolis-Hastings step to update θ , using a Gaussian random walk proposal with standard deviation $\sigma = 0.2$ and target density $p(\theta \mid z_1, x_{1:T}, y_{1:T})$. The

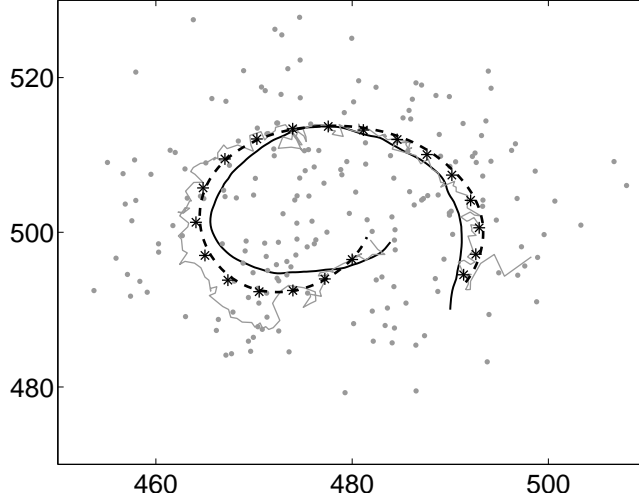


Figure 4: Target trajectory in the horizontal plane (black solid line) and smoothed posterior means for PG-AS (dashed line) and PMMH (*) under parameter uncertainty. The gray line is the PF estimate and the dots show the range-bearing measurements, transformed to Cartesian coordinates.

initial state of the system is unknown, so the variable z_1 is seen as a part of the system state. That is, the SMC sampler targets the sequence of densities $p_\theta(z_1, x_{1:t} | y_{1:t})$ for $t = 1, \dots, T$.

It is worth to point out that this SMC sampler is not more complicated to implement than a sampler targeting the original model (16). In fact, a natural way to do the implementation is to run the sampler as if targeting $[x_t^\top \ z_t^\top]^\top$ jointly³. The difference is that the z_t -particles are seen as conditional sufficient statistics for the z_t -state (which is conditionally deterministic), similarly to how one propagates the sufficient statistics for the conditionally linear state in an RBPF. The difference lies in how the backward sampling is done, where in the marginal model we only consider the x_t -states when computing the backward weights.

The PG-AS sampler was run with $N = 5$ particles for 50000 iterations, with the first 10000 samples discarded as burnin. We used an adaptive truncation level with $\gamma = 0.1$ and $\tau = 10^{-2}$ (same as before), resulting in an average truncation level of 2.3. As a comparison, we also employ a particle marginal Metropolis-Hastings (PMMH) sampler [2], with $N = 5000$ particles, also running for 50000 MCMC iterations (discarding the first 10000 samples). The smoothed estimates of the target trajectory are shown in Figure 4 and the posterior density of θ is given in Figure 5. From these results we see that the PG-AS sampler provides accurate inferential results, despite the truncation of

³For the results presented here, we used a standard bootstrap PF, which is very straightforward to implement.

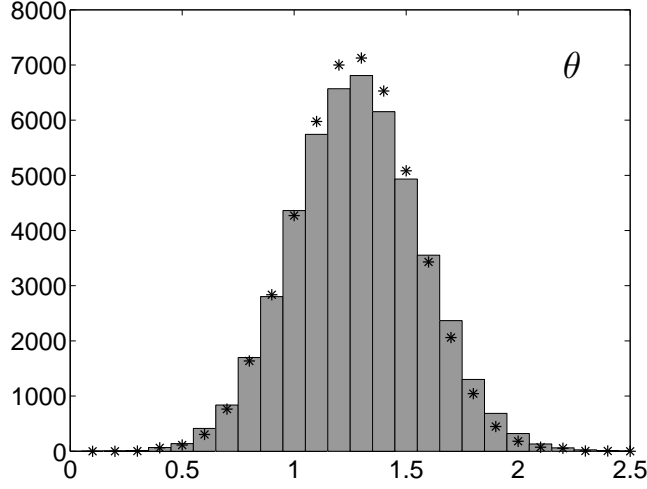


Figure 5: Histograms representing the posterior density $p(\theta \mid y_{1:T})$ for the PG-AS sampler (gray bars) and for PMMH (*). The “true” value, used in the data generation, is 1.

the backward weights and without any problem specific tuning of the variables γ and τ .

7 Discussion

PG-AS is a novel approach to PMCMC that makes use of backward simulation ideas without needing an explicit backward pass. Compared to PG-BS, a conceptually similar method that does require an explicit backward pass, PG-AS has advantages, most notably for inference in the non-Markovian SSMs that have been our focus here. When using the proposed truncation of the backward weights, we have found PG-AS to be more robust to the approximation error than PG-BS. Furthermore, for non-Markovian models, PG-AS is easier to implement than PG-BS, since it requires less bookkeeping. It can also be more memory efficient, since it does not require us to store intermediate quantities that are needed for a separate backward simulation pass, as is done in PG-BS. Finally, we note that PG-AS can be used as an alternative to PG-BS for other inference problems to which PMCMC can be applied, and we believe that it will prove attractive in problems beyond the non-Markovian SSMs that we have discussed here.

A Proof of Proposition 1

With $M = T - t$ and $w(k) = w_t^k$, the distributions of interest are given by

$$P(k) = \frac{w(k) \prod_{s=1}^M h_s(k)}{\sum_l w(l) \prod_{s=1}^M h_s(l)} \quad \text{and} \quad \tilde{P}_p(k) = \frac{w(k) \prod_{s=1}^p h_s(k)}{\sum_l w(l) \prod_{s=1}^p h_s(l)},$$

respectively. Let $\varepsilon_s \triangleq \max_{k,l} (h_s(k)/h_s(l) - 1) \leq A \exp(-cs)$ and consider

$$\begin{aligned} \left(\sum_l w(l) \prod_{s=1}^p h_s(l) \right) \prod_{s=p+1}^M h_s(k) &\leq \sum_l w(l) \prod_{s=1}^p h_s(l) \prod_{s=p+1}^M h_s(l) (1 + \varepsilon_s) \\ &= \left(\sum_l w(l) \prod_{s=1}^M h_s(l) \right) \prod_{s=p+1}^M (1 + \varepsilon_s). \end{aligned}$$

It follows that the KLD is bounded according to,

$$\begin{aligned} D_{\text{KLD}}(P \parallel \tilde{P}_p) &= \sum_k P(k) \log \frac{P(k)}{\tilde{P}_p(k)} \\ &= \sum_k P(k) \log \left(\frac{\prod_{s=p+1}^M h_s(k) (\sum_l w(l) \prod_{s=1}^p h_s(l))}{\sum_l w(l) \prod_{s=1}^M h_s(l)} \right) \\ &\leq \sum_k P(k) \sum_{s=p+1}^M \log(1 + \varepsilon_s) \leq \sum_{s=p+1}^M \varepsilon_s \leq A \sum_{s=p+1}^M \exp(-cs) \\ &= A \frac{e^{-c(p+1)} - e^{-c(M+1)}}{1 - e^{-c}}. \end{aligned} \quad \square$$

References

- [1] F. Lindsten, M. I. Jordan, and T. B. Schön, “Ancestor sampling for particle Gibbs,” in *Proceedings of the 2012 Conference on Neural Information Processing Systems (NIPS)*, Lake Tahoe, NV, USA, Dec. 2012.
- [2] C. Andrieu, A. Doucet, and R. Holenstein, “Particle Markov chain Monte Carlo methods,” *Journal of the Royal Statistical Society: Series B*, vol. 72, no. 3, pp. 269–342, 2010.
- [3] N. Whiteley, C. Andrieu, and A. Doucet, “Efficient Bayesian inference for switching state-space models using discrete particle Markov chain Monte Carlo methods,” Bristol Statistics Research Report 10:04, Tech. Rep., 2010.
- [4] F. Lindsten and T. B. Schön, “On the use of backward simulation in the particle Gibbs sampler,” in *Proceedings of the 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Kyoto, Japan, Mar. 2012.

- [5] A. Doucet and A. Johansen, “A tutorial on particle filtering and smoothing: Fifteen years later,” in *The Oxford Handbook of Nonlinear Filtering*, D. Crisan and B. Rozovsky, Eds. Oxford University Press, 2011.
- [6] M. K. Pitt and N. Shephard, “Filtering via simulation: Auxiliary particle filters,” *Journal of the American Statistical Association*, vol. 94, no. 446, pp. 590–599, 1999.
- [7] D. A. V. Dyk and T. Park, “Partially collapsed Gibbs samplers: Theory and methods,” *Journal of the American Statistical Association*, vol. 103, no. 482, pp. 790–796, 2008.
- [8] N. Whiteley, “Discussion on Particle Markov chain Monte Carlo methods,” *Journal of the Royal Statistical Society: Series B*, 72(3), p 306–307, 2010.
- [9] R. Chen and J. S. Liu, “Mixture Kalman filters,” *Journal of the Royal Statistical Society: Series B*, vol. 62, no. 3, pp. 493–508, 2000.
- [10] A. Doucet, S. J. Godsill, and C. Andrieu, “On sequential Monte Carlo sampling methods for Bayesian filtering,” *Statistics and Computing*, vol. 10, no. 3, pp. 197–208, 2000.
- [11] T. Schön, F. Gustafsson, and P.-J. Nordlund, “Marginalized particle filters for mixed linear/nonlinear state-space models,” *IEEE Transactions on Signal Processing*, vol. 53, no. 7, pp. 2279–2289, Jul. 2005.
- [12] S. Särkkä, P. Bunch, and S. Godsill, “A backward-simulation based Rao-Blackwellized particle smoother for conditionally linear Gaussian models,” in *Proceedings of the 16th IFAC Symposium on System Identification*, Brussels, Belgium, Jul. 2012.
- [13] W. Fong, S. J. Godsill, A. Doucet, and M. West, “Monte Carlo smoothing with application to audio signal enhancement,” *IEEE Transactions on Signal Processing*, vol. 50, no. 2, pp. 438–449, Feb. 2002.
- [14] N. L. Hjort, C. Holmes, P. Mller, and S. G. Walker, Eds., *Bayesian Non-parametrics*. Cambridge University Press, 2010.
- [15] S. N. MacEachern, M. Clyde, and J. S. Liu, “Sequential importance sampling for nonparametric Bayes models: The next generation,” *The Canadian Journal of Statistics*, vol. 27, no. 2, pp. 251–267, 1999.
- [16] P. Fearnhead, “Particle filters for mixture models with an unknown number of components,” *Statistics and Computing*, vol. 14, pp. 11–21, 2004.
- [17] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [18] A. Bouchard-Côté, S. Sankararaman, and M. I. Jordan, “Phylogenetic inference via sequential Monte Carlo,” *Systematic Biology*, vol. 61, no. 4, pp. 579–593, 2012.

- [19] Y. W. Teh, H. Daumé III, and D. Roy, “Bayesian agglomerative clustering with coalescents,” *Advances in Neural Information Processing*, pp. 1473–1480, 2008.
- [20] G. J. Bierman, “Fixed interval smoothing with discrete measurements,” *International Journal of Control*, vol. 18, no. 1, pp. 65–75, 1973.
- [21] S. J. Godsill, A. Doucet, and M. West, “Monte Carlo smoothing for nonlinear time series,” *Journal of the American Statistical Association*, vol. 99, no. 465, pp. 156–168, Mar. 2004.
- [22] M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, “A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking,” *IEEE Transactions on Signal Processing*, vol. 50, no. 2, pp. 174–188, 2002.
- [23] M. S. Arulampalam, B. Ristic, N. Gordon, and T. Mansell, “Bearings-only tracking of manoeuvring targets using particle filters,” *EURASIP Journal on Applied Signal Processing*, vol. 15, pp. 2351–2365, 2004.
- [24] B. Ristic, S. Arulampalam, and N. Gordon, *Beyond the Kalman filter: particle filters for tracking applications*. London, UK: Artech House, 2004.